

A conditional compression distance that unveils insights of the genomic evolution

Diogo Pratas and Armando J. Pinho

IEETA / Dept of Electronics, Telecommunications and Informatics
University of Aveiro, 3810–193 Aveiro, Portugal
pratas@ua.pt — ap@ua.pt

Abstract

We describe a compression-based distance for genomic sequences. Instead of using the usual conjoint information content, as in the classical Normalized Compression Distance (NCD), it uses the conditional information content. To compute this Normalized Conditional Compression Distance (NCCD), we need a normal conditional compressor, that we built using a mixture of static and dynamic finite-context models. Using this approach, we measured chromosomal distances between *Hominidae* primates and also between *Muroidea* (rat and mouse), observing several insights of evolution that so far have not been reported in the literature.

Introduction

The high-throughput sequencing technologies are creating an avalanche of genomic and metagenomic sequences, nonexistent a few years ago. We are now able to computationally evaluate similarities, or their absence, among species and across different regions of the same species, using whole genomes.

Common biological approaches for determining distances, usually using FISH techniques, are very expensive and time-consuming. Computational approaches have emerged as an affordable, fast and automated process to deal with this problem. Several computational distance metrics have been proposed, where some of the most popular are the Hamming [1] and Levenshtein [2] distances. The Hamming distance can only be applied when the sequences are aligned with precision and have the same size, requirements hardly found in large genomic sequences. The Levenshtein distance explores transformations between the sequences, namely insertions, deletions and substitutions. Although quite successful, its computational time is prohibitive for large sequences (the fastest known implementation runs with time complexity $O(n^2/\log n)$).

Compression-based approaches emerged as a natural way for measuring distances, because, together with the appropriate decoder, the bitstream produced by a loss-less compression algorithm allows the reconstruction of the original data and, therefore, can be seen as an upper bound of the algorithmic entropy of the sequence. A compression-based distance computes the distance between two objects using the number of bits needed to describe one of them when a description of the other is available, as well as the number of bits required to describe each of them.

Compression-based distances are founded on the Kolmogorov notion of complexity, also known as algorithmic entropy, where $K(x)$ of a string x is the length of the

shortest binary program x^* that computes x in an appropriate universal Turing machine and halts [3]. As such, $K(x) = |x^*|$, the length of x^* , denotes the number of bits of information from which x can be computationally retrieved [4]. The conditional Kolmogorov complexity, $K(x|y)$, denotes the length of the shortest binary program, in the universal prefix Turing machine, that on input y outputs x . A special case occurs when y is an empty string, $y = \lambda$, and hence $K(x|\lambda) = K(x)$.

Bennett introduced the information distance [5], $E(x, y) = \max\{K(x|y), K(y|x)\}$, defined as the length of the shortest binary program for the reference universal prefix Turing machine that with input x computes y , as well as with y computes x . The normalized version (NID [6]) of $E(x, y)$ is defined as

$$\text{NID}(x, y) = \frac{\max\{K(x|y), K(y|x)\}}{\max\{K(x), K(y)\}}, \quad (1)$$

up to an additive logarithmic term. The normalized compression distance (NCD) [7] emerged to efficiently compute the NID, due to the non-computability of K ,

$$\text{NCD}(x, y) = \frac{C(x, y) - \min\{C(x), C(y)\}}{\max\{C(x), C(y)\}}, \quad (2)$$

up to an additive logarithmic term, where $C(x)$ and $C(y)$ represent, respectively, the number of bits of the compressed version of x and y , and $C(x, y)$ the number of bits of the conjoint compression of x and y (usually, x and y are concatenated). Distances near one indicate dissimilarity, while distances near zero indicate similarity. It can be seen that for $\text{NCD}(x, x) = 0$ to hold, then the compressor needs to verify $C(x, x) \approx C(x)$, one of the most important properties of a *normal* compressor [7].

In this paper, we describe an admissible normalized compression distance, relying on a conditional compressor, that builds an internal model of the data using a mixture of static and dynamic finite-context models (FCMs). We assess the metric and its inherent parameterized compressor, and we present some results of chromosomal distances between several large eukaryotic chromosomes, namely *Hominidae* primates and *Muroidea*, confirming several documented results and pointing out some undocumented observations.

Proposed Approach

A direct substitution of K by C in (1) would require the availability of compressors that are able to produce conditional compression, i.e., $C(x|y)$ and $C(y|x)$. Most compressors do not have this functionality and, therefore, the NCD avoids it by using suitable manipulations of (1) [7]. Instead of $C(x|y)$ and $C(y|x)$, a term corresponding to the conjoint compression of x and y , $C(x, y)$, was preferred. Usually, this $C(x, y)$ term is interpreted as the compression of the concatenation of x and y , but, in fact, it could be any other form of combination between x and y . Concatenation is often used because it is easy to obtain, but in fact its use may hamper the efficiency of the measure [8].

To overcome this limitation, we propose use the direct form, to which we call the Normalized Conditional Compression Distance (NCCD),

$$\text{NCCD}(x, y) = \frac{\max\{C(x|y), C(y|x)\}}{\max\{C(x), C(y)\}}, \quad (3)$$

where “Conditional” means that the compressor C needs to be able to perform conditional compression.

The conditional compressor

We have built a NCCD compressor based on two model classes (we call them “static” and “dynamic”), each one composed of mixtures of finite-context models (FCMs) of several orders [9–11]. To compute $C(x|y)$, the compression is performed in two phases. In the first phase, the static class of FCMs accumulates the counts regarding the y object. After the entire y object was processed, the models are kept frozen and, hence, the second phase starts. At this point, the x object starts to be compressed using the static models computed during the first phase, in cooperation with the set of FCMs of the dynamic class, that dynamically accumulate the counts only from x .

The probability of each symbol is obtained by mixing the probabilities provided by each FCM of the static and dynamic models, using a weighted average, according to

$$P(x_{n+1}) = \sum_k P(x_{n+1}|x_{n-k+1..n}) w_{k,n}, \quad (4)$$

where $w_{k,n}$ denotes the weight assigned to the finite-context model k and $\sum_k w_{k,n} = 1$. The conditional probabilities are given by the estimator

$$P(s|x_{n-k+1..n}) = \frac{C(s|x_{n-k+1..n}) + \alpha}{C(x_{n-k+1..n}) + 4\alpha}, \quad (5)$$

where $C(s|x_{n-k+1..n})$ represents the number of times that, in the past, symbol s was found having $x_{n-k+1..n}$ as the conditioning context and where $C(x_{n-k+1..n})$ is the total number of events that has occurred so far in association with context $x_{n-k+1..n}$.

For stationary sources, we could compute weights such that $w_{k,n} = P(k|x_{1..n})$, i.e., according to the probability that model k has generated the sequence until that point. In that case, we would get

$$w_{k,n} = P(k|x_{1..n}) \propto P(x_{1..n}|k)P(k), \quad (6)$$

where $P(x_{1..n}|k)$ denotes the likelihood of sequence $x_{1..n}$ being generated by model k and $P(k)$ denotes the prior probability of model k . Assuming $P(k) = 1/K$, where K denotes the total number of FCMs, we obtain $w_{k,n} \propto P(x_{1..n}|k)$. Calculating the logarithm we get

$$\log_2 P(x_{1..n}|k) = \log_2 \prod_{i=1}^n P(x_i|k, x_{1..i-1}) = \sum_{i=1}^n \log_2 P(x_i|k, x_{1..i-1}), \quad (7)$$

which corresponds to the code length that would be required by model k for representing the sequence $x_{1..n}$. It is, therefore, the accumulated measure of the performance of model k until instant n . However, since the DNA sequences are not stationary, a good performance of a model in a certain region of the sequence might not be attained in other regions. Hence, the performance of the models have to be measured in the recent past of the sequence, for example using a mechanism of progressive forgetting of past measures. For that, we use the recursive relation

$$\sum_{i=1}^n \log_2 P(x_i|k, x_{1..i-1}) = \quad (8a)$$

$$= \gamma \sum_{i=1}^{n-1} \log_2 P(x_i|k, x_{1..i-1}) + \log_2 P(x_n|k, x_{1..n-1}). \quad (8b)$$

This relation corresponds to a first-order recursive filter that, for $\gamma \in [0, 1)$, has a low-pass characteristic and an exponentially decaying impulse response. For additional information, see, for example, [12, 13].

Parameterization and assessment

The parameters used in each compression measure must be kept constant, in order to be used as a valid comparable metric between distances (otherwise it will change the meaning of C). Accordingly, we have used a fixed setup of five static FCMs and three dynamic FCMs, mixed using a set of weights estimated with $\gamma = 0.9$. From our experience, we have verified that $\gamma = 0.99$ maximizes the compression gain for bacterial genomes, while for eukaryotic genomes $\gamma = 0.9$ seems to be the best choice. The orders used for the static models were: 4, 6, 8, 10 and 15. For the first four we used $\alpha = 1$ (Laplace estimator), whereas the one with the highest order we used $\alpha = 0.001$. Usually, a small α is important only for high orders (above ten). Moreover, the high order used (15) ensures an admissible identity (i.e., $\text{NCCD}(x, x) \approx 0$), as Fig. 1 suggests. The curve in Fig. 1 labeled “*lossy*” corresponds to using always the best FCM for each base and shows that the first part of the curves is due to the adaptation of the method when not enough data is present, suggesting that a very small sequence may harm the identity property, also observed in very large sequences. The latter drawback may be overcome using higher FCM orders, at the cost of additional computational memory.

The three FCMs of the dynamic class have orders 4, 10 and 15, where the first two rely on a Laplace probability estimator and the last one use $\alpha = 0.05$. For the two deeper models, the inverted repeats are also taken into account [14]. The maximum counters used in each static model were, respectively, 2^9 , 2^{12} , 2^{12} . This limitation acts also as a forgetting mechanism, because the counters are divided by two when one of them reaches the maximum, decreasing the importance of statistics collected in the far past. More information regarding FCM parameterization can be obtained in [12, 13, 15].

The DNA data sequences are products of sequencing techniques, which have a sequencing quality, coverage and assembly technique associated [16]. Although these

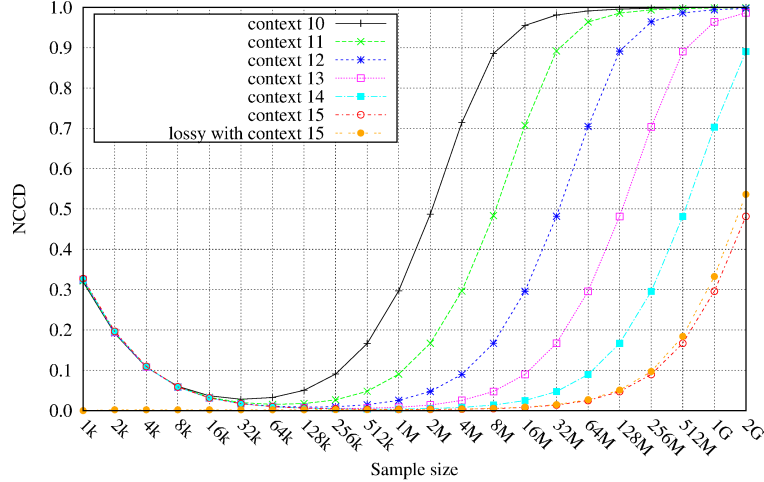


Figure 1: $NCCD(x, x)$ value on uniformly distributed DNA (synthetic) sequences with custom sizes, for several depths of the highest order model. The “lossy” curve shows the behavior of NCCD when the best FCM is chosen for each base, corresponding to a lower bound of the (non-reversible) compressor.

external factors may sometimes constitute a problem, we believe that generally they are mitigated by the compressor [17]. Nevertheless, since we use a metric based on conditionals targeting genomic sequences, we have assessed the impact of uniformly distributed mutations, namely substitutions, insertions and deletions, over 50 MB of *real* (first 50 MB of chromosome 1 from *H. sapiens*) and synthetic (simulated using XS from Exon [18]) genomic data, as can be seen in the top graph of Fig. 2. Substitutions seem to be the most difficult mutation type to be handled by the compressor, although only slightly, and, hence, by the NCCD. Although it is clear that the method is still reporting reasonable distances for sequences with 10% of mutations, both for the *real* and synthetic sequences.

Finally, we have assessed the importance of sequence completeness using progressive missing data, as the bottom graph of Fig. 2 depicts. As expected, it is characterized by an approximately linear behavior. However, there is a gap between the curves of the *real* and synthetic sequences, specially when there are lower missing rates. This is due to the nature of the sequences, namely the self-similarity, since the beginning of the *real* sequence is composed by a telomeric zone (highly-repetitive). On the other hand, the synthetic sequence does not yield an exact zero of the NCCD when the missing rate is zero, because it has been simulated with several approximately repeating zones. This may be overcome with higher FCM orders, although at the cost of more computer memory.

Results

The data set is composed of six genomes (Table 1), downloaded from the NCBI website (<ftp://ftp.ncbi.nlm.nih.gov/genomes>).

Figure 3 presents the inter-chromosomal NCCD distance heatmaps relatively to *H. sapiens* with the rest of the primates and *M. musculus*, and *M. musculus* relatively

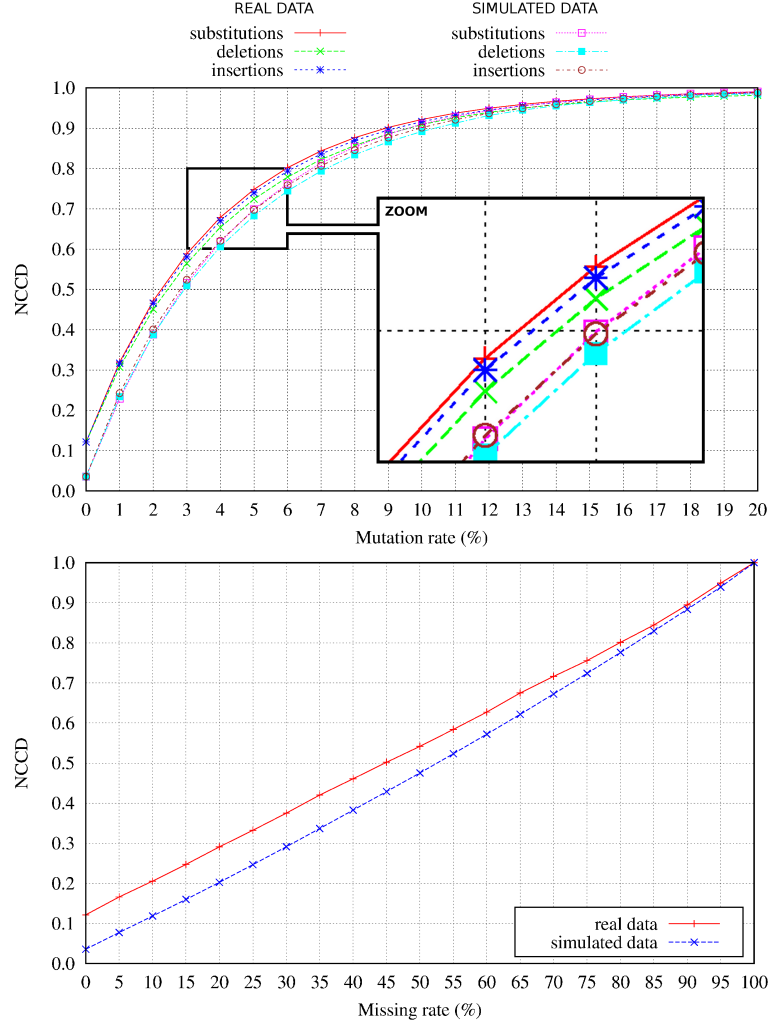


Figure 2: NCCD performance on synthetic and *real* 50 MB of genomic mutated data (top) and on progressive block missing data (bottom).

to *R. norvegicus*, plotted in an *all with all* scheme. As can be seen, for all primate species there is a direct correlation with the respective chromosomal number, with the exception of chromosome 2 (related to 2A and 2B). This is justified by a presumed chromosomal fusion in humans from previous ancestors [19].

Moreover, the human Y chromosome is highly related with the X chromosome of other primate species, namely the *P. troglodytes*, because the Y chromosome exchanged genetic information with X in the recombination process [20]. Furthermore, there is a low distance between chromosomes 5 and 17 of the *G. gorilla* and *H. sapiens*, justified by a chromosomal translocation [21].

Relatively to *M. musculus*, there is an obvious similarity with *R. norvegicus*, although smaller than in *P. maniculatus* / *M. norvegicus* [22]. When compared with the primates, no important similarities are found (at a genomic level), specially in human chromosomes 19 and 22. Moreover, it seems that only the mitochondrial sequences attain some level of similarity. Nevertheless, the *M. musculus* (MM) and

Table 1: Data set used in the experiments. The number of expected chromosome pairs for each species is represented by “Exp”, while “Missing” is a nonexistence sequence and Mb represents the approximated size in Mega bases.

Organism	Build	Exp	Missing	Mb
<i>Homo sapiens</i>	37.p10	23	-	2,861
<i>Pan troglodytes</i>	2.1.4	24	-	2,756
<i>Gorilla gorilla</i>	r100	24	Y	2,719
<i>Pongo abelii</i>	1.3	24	Y	3,028
<i>Mus musculus</i>	38.p1	20	-	2,716
<i>Rattus norvegicus</i>	5.1	21	Y	2,443

R. norvegicus (RN) *diagonal* is very dissipated for such a low distance depicted in the mitochondrial sequence. In fact, only chromosomes (C) 18 and X seem to be homologous (in the *diagonal*). Subsequent analysis show strong similarity between MM C2 / RN C3, MM C9 / RN C8 and MM C11 / RN C10, and considerable similarity between MM C4 / RN C5, MM C6 / RN C4, MM C12 / RN C6 and MM C14 / RN C15, without detracting other important patterns.

Figure 4 presents the chromosomal distances of *P. troglodytes*, *G. gorilla* and *P. abelii* (chromosomes 2A and 2B have been concatenated) according to the *H. sapiens* chromosomes order. At glance, *P. troglodytes* has the lowest distance relatively to *H. sapiens*, followed by *G. gorilla* and *P. abelii*, respectively. Specifically, the *G. gorilla* chromosomes 5 and 17 have large distances because of the previous mentioned translocation, while *P. abelii* seems to have a very different chromosome 1, besides other relevant dissimilarities.

According to [23], besides the high divergence of Y chromosome, there are several breakpoints in chromosomes 4, 5 and 12, which were tested by fluorescence *in situ* hybridization (FISH) in *P. troglodytes*, using *H. sapiens* as reference. Figure 4 reports the same dissimilarities, surprisingly adding chromosome 17.

Finally, we have found that chromosomes 4, 12 and 18 of *G. gorilla* have lower distances to *H. sapiens* than to the respective *P. troglodytes* chromosomes, while chromosomes 5 and 17 of *G. gorilla* have higher distances than those of *P. abelii*. Mitochondrial sequences, as expected, show that *P. troglodytes* is the nearest species to *H. sapiens*, followed by the *G. gorilla* and *P. abelii*.

Conclusion

We have described a compressed-based metric for measuring distances between genomic sequences, based on the conditional information content. This approach requires a *normal* conditional compressor, that we have defined and assessed in this work. The compressor is constituted by a set of multiple static and dynamic finite-context models, that cooperate under a supervision mixture model. It is able to handle several types of mutations, and hence rendering it a good candidate to study large eukaryotic chromosomes. We have calculated chromosomal distances between

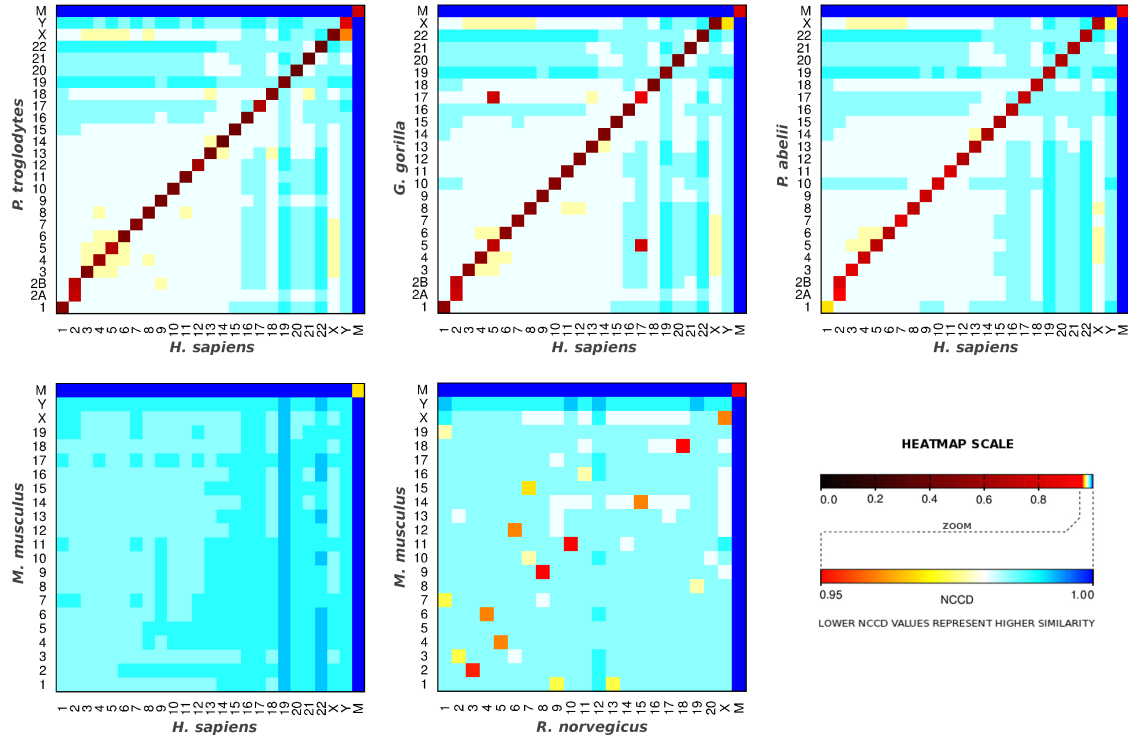


Figure 3: *P. troglodytes*, *G. gorilla*, *P. abelii* and *M. musculus* inter-genomics chromosomal NCCD heatmaps, in relation to *H. sapiens*, and *M. musculus* in relation to *R. norvegicus*.

Hominidae primates and also *Muroidea* (rat and mouse) rodents superfamily, attaining results that agree with several already documented results, mainly using expensive and time-consuming FISH approaches, but also unveiling undocumented ones.

Acknowledgements

This work was supported in part by FEDER through the Operational Program Competitiveness Factors - COMPETE and by National Funds through FCT - Foundation for Science and Technology, in the context of the projects FCOMP-01-0124-FEDER-022682 (FCT reference PEst-C/EEI/UI0127/2011) and Incentivo/EEI/UI0127/2013.

References

- [1] R. Hamming, "Error detecting and error correcting codes," *Bell System Technical Journal*, vol. 29, no. 2, pp. 147–160, 1950.
- [2] V. Levenshtein, "Binary codes capable of correcting deletions, insertions and reversals," in *Soviet physics doklady*, vol. 10, 1966, p. 707.
- [3] A. Turing, "On computable numbers, with an application to the Entscheidungsproblem," *Proceedings of the London Mathematical Society*, vol. 42, no. 2, pp. 230–265, 1936.
- [4] M. Li and P. Vitányi, *An introduction to Kolmogorov complexity and its applications*, 3rd ed. Springer, 2008.

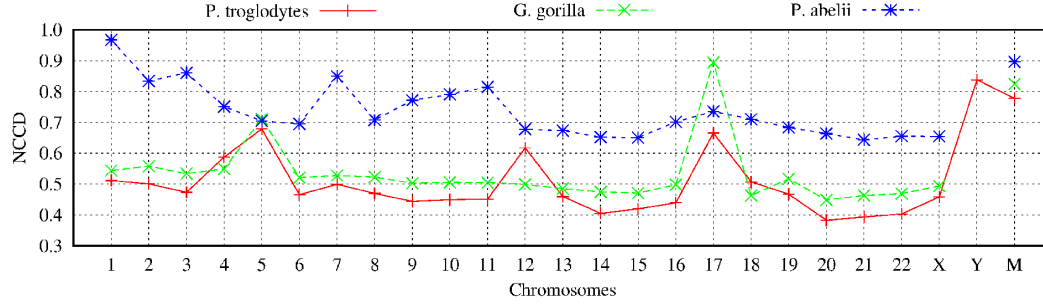


Figure 4: *P. troglodytes*, *G. gorilla* and *P. abelii* related chromosomal NCCD values using *H. sapiens* as reference.

- [5] C. H. Bennett, P. Gács, M. L. P. M. B. Vitányi, and W. H. Zurek, "Information distance," *IEEE Trans. on Information Theory*, vol. 44, no. 4, pp. 1407–1423, Jul. 1998.
- [6] M. Li, X. Chen, X. Li, B. Ma, and P. M. B. Vitányi, "The similarity metric," *IEEE Trans. on Information Theory*, vol. 50, no. 12, pp. 3250–3264, Dec. 2004.
- [7] R. Cilibrasi and P. M. B. Vitányi, "Clustering by compression," *IEEE Trans. on Information Theory*, vol. 51, no. 4, pp. 1523–1545, Apr. 2005.
- [8] M. Cebrián, M. Alfonseca, and A. Ortega, "Common pitfalls using the normalized compression distance: what to watch out for in a compressor," *Communications in Information and Systems*, vol. 5, no. 4, pp. 367–384, 2005.
- [9] T. C. Bell, J. G. Cleary, and I. H. Witten, *Text compression*. Prentice Hall, 1990.
- [10] D. Salomon, *Data compression - The complete reference*, 4th ed. Springer, 2007.
- [11] K. Sayood, *Introduction to data compression*, 4th ed. Morgan Kaufmann, 2012.
- [12] D. Pratas and A. J. Pinho, "Compressing the human genome using exclusively Markov models," in *Advances in Intelligent and Soft Computing, Proc. of the 5th Int. Conf. on Practical Applications of Computational Biology & Bioinformatics, PACBB 2011*, vol. 93, Apr. 2011, pp. 213–220.
- [13] A. J. Pinho, D. Pratas, and P. J. S. G. Ferreira, "Bacteria DNA sequence compression using a mixture of finite-context models," in *Proc. of the IEEE Workshop on Statistical Signal Processing*, Nice, France, Jun. 2011.
- [14] A. J. Pinho, A. J. R. Neves, and P. J. S. G. Ferreira, "Inverted-repeats-aware finite-context models for DNA coding," in *Proc. of the 16th European Signal Processing Conf., EUSIPCO-2008*, Lausanne, Switzerland, Aug. 2008.
- [15] A. J. Pinho, P. J. S. G. Ferreira, A. J. R. Neves, and C. A. C. Bastos, "On the representability of complete genomes by multiple competing finite-context (Markov) models," *PLoS ONE*, vol. 6, no. 6, p. e21588, 2011.
- [16] D. Church, M. Deanna, V. Schneider *et al.*, "Modernizing reference genome assemblies," *PLoS Biology*, vol. 9, no. 7, p. e1001091, 2011.
- [17] M. Cebrián, M. Alfonseca, and A. Ortega, "The normalized compression distance is resistant to noise," *IEEE Trans. on Information Theory*, vol. 53, no. 5, pp. 1895–1900, 2007.
- [18] D. Pratas, A. J. Pinho, and S. Garcia, "Exon: A web-based software toolkit for dna sequence analysis," in *6th International Conference on Practical Applications of Computational Biology & Bioinformatics*. Springer, 2012, pp. 217–224.

- [19] J. Ijdo, A. Baldini, D. Ward, S. Reeders, and R. Wells, "Origin of human chromosome 2: an ancestral telomere-telomere fusion." *Proceedings of the National Academy of Sciences USA*, vol. 88, no. 20, pp. 9051–9055, 1991.
- [20] J. Hughes *et al.*, "Chimpanzee and human Y chromosomes are remarkably divergent in structure and gene content," *Nature*, vol. 463, no. 7280, pp. 536–539, 2010.
- [21] R. Samonte and E. Eichler, "Segmental duplications and the evolution of the primate genome," *Nature Reviews Genetics*, vol. 3, no. 1, pp. 65–72, 2002.
- [22] C. Ramsdell *et al.*, "Comparative genome mapping of the deer mouse (*Peromyscus maniculatus*) reveals greater similarity to rat (*Rattus norvegicus*) than to the lab mouse (*Mus musculus*)," *BMC Evolutionary Biology*, vol. 8, no. 1, p. 65, 2008.
- [23] T. Mikkelsen *et al.*, "Initial sequence of the chimpanzee genome and comparison with the human genome." *Nature*, 2005.